

A. Pope

PATENT

Atty. Docket No. SVH-003

(7453/6)

Box Patent Application

Assistant Commissioner for Patents

Washington, D.C. 20231

NEW APPLICATION TRANSMITTAL

Transmitted herewith for filing is the patent application of

Inventors: Jamie Callan

WARNING: Patent must be applied for in the name(s) of all of the actual inventor(s). 37 CFR 1.41(a) and 1.53(b)

For (title): System and Method For Filtering A Document Stream

1. Type of Application

This new application is for a(n) (check one applicable item below):

☒ Original

☐ Design

☐ Plant

WARNING: Do not use this transmittal for a completion in the U.S. of an International Application under 35 U.S.C. 371(c)(4) unless the International Application is being filed as a divisional, continuation or continuation-in-part application.

NOTE: If one of the following 3 items apply then complete and attach ADDED PAGES FOR NEW APPLICATION TRANSMITTAL WHERE BENEFIT OF A PRIOR U.S. APPLICATION CLAIMED.

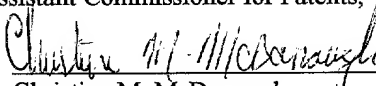
☐ Divisional

☐ Continuation

☐ Continuation-in-Part

CERTIFICATE OF EXPRESS MAILING UNDER 37 C.F.R. 1.10

I hereby certify that the attached document, and any documents referred to as enclosed herein, are being deposited with the United States Postal Service, postage prepaid, on August 18, 1997 utilizing the "Express Mail Post Office to Addressee" service of the United States Postal Service, mailing label number EM354713966US, in an envelope addressed to Box Patent Application, Assistant Commissioner for Patents, Washington, DC 20231.


Christine M. McDonough

Express Mail Label No. EM354713966US

2. Benefit of Prior U.S. Application(s) (35 USC 120)

NOTE: If the new application being transmitted is a divisional, continuation or a continuation-in-part of a parent case, or where the parent case is an International Application which designated the U.S., then check either the first option below or the second option below. If the second option is checked, the ADDED PAGES FOR NEW APPLICATION TRANSMITTAL WHERE BENEFIT OF PRIOR U.S. APPLICATION (S) IS CLAIMED must be completed and attached.

- ☐ The new application transmitted claims the benefit of prior U.S. application(s) and the priority information is contained in the enclosed new application
- ☐ The new application being transmitted claims the benefit of prior U.S. application(s) and enclosed are ADDED PAGES FOR NEW APPLICATION TRANSMITTAL WHERE BENEFIT OF PRIOR U.S. APPLICATION(S) CLAIMED.
- ☐ Amend the Specification by inserting before the first line the sentence:

"This is a

- ☐ continuation
- ☐ continuation-in-part
- ☐ divisional

of copending application(s)

- ☐ serial number 0 / filed on "
- ☐ International Application filed on and which designated the U.S."

3. Priority Claimed on Provisional Application(s) Under 35 U.S.C. 119(e).

- ☐ The new application claims benefit of priority on the following U.S. Provisional Application(s):

Application No.

Filed

4. Priority Claimed on Prior International Application(s) Under 35 U.S.C. 119.

Country

Application No.

Filed

5. Papers Enclosed Which Are Required For Filing Date Under 37 CFR 1.53(b) (Regular) or 37 CFR 1.153 (Design) Application

17 Pages of specification

3 Pages of claims

1 Pages of Abstract

2 Sheets of drawing

New Application Transmittal

Page 3

- ☐ formal
- ☒ informal

WARNING: DO NOT submit original drawings. A high quality copy of the drawings should be supplied when filing a patent application. The drawings that are submitted to the Office must be on strong, white, smooth, and non-shiny paper and meet the standards according to § 1.84. If corrections to the drawings are necessary, they should be made to the original drawing and a high-quality copy of the corrected original drawing then submitted to the Office. **Only one copy is required or desired.** Comments on proposed new 37 CFR 1.84. Notice of March 9, 1988 (1990 O.G. 57-62).

NOTE: "Identifying indicia such as the serial number, group and unit, title of the invention, attorney's docket number, inventor's name, number of sheets, etc., not to exceed 2-3/4 inches (7.0 cm.) in width may be placed in a centered location between the side edges within three fourths inch (19.1 mm.) of the top edge. Either this marking technique on the front of the drawing or the placement, although not preferred, of this information and the title of the invention on the back of the drawings is acceptable." Proposed 37 CFR 1.84(1). Notice of March 9, 1988 (1090 O.G. 67-62).

6. Additional papers enclosed

- ☐ Preliminary Amendment
- ☐ Information Disclosure Statement
- ☐ Form PTO-1449
- ☐ Citations
- ☐ Declaration of Biological Deposit
- ☐ Submission of "Sequence Listing," computer readable copy and/or amendment pertaining thereto for biotechnology invention containing nucleotide and/or amino acid sequence.
- ☐ Authorization of Attorney(s) to Accept and Follow Instructions from Representative
- ☐ Special Comments
- ☐ Other

7. Declaration or oath

- ☒ Enclosed but unexecuted.
- ☐ Enclosed
- executed by (check all applicable boxes)

- ☐ inventor(s).
- ☐ legal representative of inventor(s). 37 CFR 1.42 or 1.43
- ☐ joint inventor or person showing a proprietary interest on behalf of inventor who refused to sign or cannot be reached.

☐ this is the petition required by 37 CFR 1.47 and the statement required by 37 CFR 1.47 is also attached. See item 13 below for fee.

- ☐ Not Enclosed.

New Application Transmittal

Page 4

WARNING:

Where the filing is a completion in the U.S. of an International Application but where a declaration is not available or where the completion of the U.S. application contains subject matter in addition to the International Application the application may be treated as a continuation or continuation-in-part, as the case may be, utilizing ADDED PAGE FOR NEW APPLICATION TRANSMITTAL WHERE BENEFIT OF PRIOR U.S. APPLICATION CLAIMED.

- ☐ Application is made by a person authorized under 37 CFR 1.41(c) on behalf of all the above named inventor(s). The declaration or oath, along with the surcharge required by 37 CFR 1.16(e) can be filed subsequently.

NOTE: It is important that all the correct inventor(s) are named for filing under 37 CFR 1.41(c) and 1.53(b).

- ☐ Showing that the filing is authorized.
(Not required unless called into question. 37 CFR 1.41(d).

8. Inventorship Statement

WARNING: If the named inventors are each not the inventors of all the claims an explanation, including the ownership of the various claims at the time the last claimed invention was made, should be submitted.

The inventorship for all the claims in this application are:

- ☐ The same
- ☐ Are not the same. An explanation, including the ownership of the various claims at the time the last claimed invention was made,
- ☐ is submitted.
- ☐ will be submitted.

9. Language

NOTE: An application including a signed oath or declaration may be filed in a language other than English. A verified English translation of the non-English language application and the processing fee of \$130.00 required by 37 CFR 1.17(k) is required to be filed with the application or within such time as may be set by the Office. 37 CFR 1.52(d).

NOTE: A non-English oath or declaration in the form provided or approved by the PTO need not be translated. 37 CFR 1.69(b).

- ☒ English
- ☐ non-English
- ☐ the attached translation is a verified translation. 37 CFR 1.52(d).

10. Assignment

- ☒ An assignment of the invention to Sovereign Hill Software, Inc.

New Application Transmittal

Page 5

☐ is (are) attached. A separate "ASSIGNMENT COVER LETTER ACCOMPANYING NEW PATENT APPLICATION" is also attached.

☒ will follow.

NOTE: "If an assignment is submitted with a new application, send two separate letters -- one for the application and one for the assignment." Notice of May 4, 1990 (1114 D.G. 77-78).

11. Certified Copy Certified copy(ies) of the application(s)

Country	Application No.	Filed
Country	Application No.	Filed
Country	Application No.	Filed
Country	Application No.	Filed
Country	Application No.	Filed
Country	Application No.	Filed

from which priority is claimed

☐ is (are) attached.

☐ will follow.

NOTE: The foreign application forming the basis for the claim for priority **must** be referred to in the **oath** or **declaration**. 37 CFR 1.55(a) and 1.63.

NOTE: This item is for any foreign priority for which the application being filed directly relates. If any parent U.S. application or International Application from which this application claims benefit under 35 U.S.C. 120 is itself entitled to priority from a prior foreign application then complete item 18 on the ADDED PAGES FOR NEW APPLICATION TRANSMITTAL WHERE BENEFIT OF PRIOR U.S. APPLICATION(S) CLAIMED.

12. Fee Calculation (37 CFR 1.16)

A. ☒ Regular application

CLAIMS AS FILED

	Number Filed	Number Extra		Rate	Basic Fee 37 CFR 1.16(a) \$770.00
Total Claims (37 CFR 1.16 (c))	16	- 20 = 0	X	\$ 22.00	\$
Independent Claims (37 CFR 1.16 (b))	3	- 3 = 1	X	\$ 80.00	\$
Multiple Dependent Claim(s), If any (37 CFR 1.16(d))			+	\$ 260.00	\$

New Application Transmittal

Page 6

- ☐ Amendment canceling extra claims enclosed.
- ☐ Amendment deleting multiple-dependencies enclosed.
- ☐ Fee for extra claims is not being paid at this time.

NOTE: If the fees for extra claims are not paid on filing they must be paid or the claims canceled by amendment, prior to the expiration of the time period set for response by the Patent and Trademark Office in any notice of fee deficiency. 37 CFR 1.16(d).

Filing Fee Calculation \$ 770.00

- B. ☐ **Design application**
(\$320.00--37 CFR 1.16(f))

Filing Fee Calculation \$

- C. ☐ **Plant application**
(\$530.00--37 CFR 1.16(g))

Filing Fee Calculation \$

13. Small Entity Statement(s)

- ☐ Verified Statements that this is a filing by a small entity under 37 CFR 1.9 and 1.27

Filing Fee Calculation (50% of A, B or C above) \$

NOTE: Any excess of the full fee paid will be refunded if a verified statement and a refund request are filed within 2 months of the date of timely payment of a full fee. 37 CFR 1.28(a).

14. Request for International-Type Search (37 CFR 1.104(d)) (complete, if applicable)

- ☐ Please prepare an international-type search report for this application at the time when national examination on the merits takes place.

15. Fee Payment Being Made At This Time

- ☒ Not Enclosed
- ☒ No filing fee is to be paid at this time. (This and the surcharge required by 37 CFR 1.16(e) can be paid subsequently.)
- ☐ Enclosed
- ☐ basic filing fee \$
- ☐ recording assignment
(\$40.00; 37 CFR 1.21(h)) \$
- ☐ petition fee for filing by other than all the
inventors or person on behalf of the
inventor where inventor refused to sign
or cannot be reached.
(\$130.00; 37 CFR 1.47 and 1.17(h)) \$

Total fees enclosed \$

16. Method of Payment of Fees

NOTE: Fees should be itemized in such a manner that it is clear for which purpose the fees are paid. 37 CFR 1.22(b).

17. Authorization to Charge Additional Fees

WARNING: If no fees are to be paid on filing the following items should not be completed.

WARNING: Accurately count claims, especially multiple dependent claims, to avoid unexpected high charges, if extra claim charges are authorized.

NOTE: Because additional fees for excess or multiple dependent claims not paid on filing or on later presentation must only be paid or these claims canceled by amendment prior to the expiration of the time period set for response by the PTO in any notice of fee deficiency (37 CFR 1.16(d), it might be best not to authorize the PTO to charge additional claim fees, except possibly when dealing with amendments after final action.

- ☐ 37 CFR 1.16(e) (surcharge for filing the basic filing fee and/or declaration on a date later than the filing date of the application)

☐ 37 CFR 1.17 (application processing fees)

WARNING: While 37 CFR 1.17(a),(b), (c) and (d) deal with extensions of time under § 1.136(a) this authorization should be made only with the knowledge that: "Submission of the appropriate extension fee under 37 C.F.R. 1.136(a) is to no avail unless a request or petition for extension is filed." (Emphasis added). Notice of November 5, 1985 (1060 O.G. 27).

☐ 37 CFR 1.18 (issue fee at or before mailing of Notice of Allowance, pursuant to 37 CFR 1.311(b))

NOTE: Where an authorization to charge the issue fee to a deposit account has been filed before the mailing of a Notice of Allowance, the issue fee will be automatically charged to the deposit account at the time of mailing the notice of allowance, 37 CFR 1.31(b).

NOTE: 37 CFR 1.28(b) requires "Notification of any change in loss of entitlement to small entity status must be filed in the application . . . prior to paying, or at the time of paying, . . . issue fee". From the wording of 37 CFR 1.28(b): (a) notification of change of status must be made even if the fee is paid as "other than a small entity" and (b) no notification is required if the change is to another small entity.

18. Instructions As To Overpayment

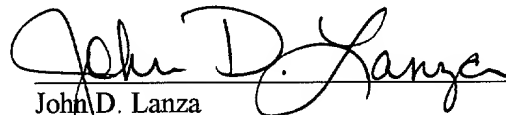
☒ credit Account No. 20-0531

☐ refund

Date: August 18, 1997

Reg. No. 40,060

Tel. No. (617) 248-7604



John D. Lanza

Attorney for Applicants

Testa, Hurwitz, & Thibault, LLP

High Street Tower

125 High Street

Boston, Massachusetts 02110

☐ **Incorporation by reference of added pages**

Check the following item if the application in this transmittal claims the benefit of prior U.S. application(s) (including an international application entering the U.S. stage as a continuation, divisional or C-I-P application) and complete and attach the ADDED PAGES FOR NEW APPLICATION TRANSMITTAL WHERE BENEFIT OF PRIOR U.S. APPLICATION(S) CLAIMED

☐ Plus Added Pages For New Application Transmittal Where Benefit of prior U.S. Application(s) Claimed

Number of pages added

☐ Plus Added Pages For Papers Referred To In Item 4 Above

Number of pages added

☐ Plus "Assignment Cover Letter Accompanying New Application"

Number of pages added

☐ **Statement Where No Further Pages Added**

(if no further pages form a part of this Transmittal then end this Transmittal with this page and check the following item)

☒ This transmittal ends with this page.

APPLICATION
FOR
UNITED STATES
LETTERS PATENT

SPECIFICATION

(For Attorney Docket No. SVH-003)

TO ALL WHOM IT MAY CONCERN:

Be it known that I, **Jamie Callan**, a citizen of the United States of America and residing at 2 Wellington Circle, Easthampton, Massachusetts 01027, have invented new and useful improvements in

SYSTEM AND METHOD FOR FILTERING A DOCUMENT STREAM

of which the following is a specification.

412JDL7453/6.397662-1

Express Mail Label No. EM354713966US

SYSTEM AND METHOD FOR FILTERING A DOCUMENT STREAM

Field of the Invention

The present invention relates to information retrieval techniques and, in particular, to techniques for efficiently filtering a document stream based on the minimum number of query terms a document must possess to satisfy a query profile.

Background of the Invention

Although statistical retrieval models are now accepted widely, there has been little research on how to adapt them to the demands of high-speed document filtering. The problems of document retrieval and document filtering are similar at an abstract level, but the architectures required, the optimizations that are possible, and the quality of the information available, are all different.

Retrieval of documents from an archival collection (retrospective retrieval) and filtering documents from an incoming stream of documents (document filtering or selective dissemination of information) have been described as two sides of the same coin. Both tasks consist of determining quickly how well a document matches an information need. Many of the underlying issues are the same; for example, deciding how to represent each document, how to describe the information need in a query language, what words to ignore (e.g., stop words), whether or not to stem words, and how to interpret evidence of relevance.

Many document filtering techniques are based on the assumption that effective document retrieval techniques are also effective document filtering techniques. However, when filtering research is conducted with a retrieval system,

important issues can be overlooked. Different architectures are possible, and perhaps required, to rapidly compare persistent information needs to transient documents. A filtering algorithm must make decisions based upon incomplete information; it may know what has happened in the past, but it generally cannot
5 know, nor can it generally wait to know, what documents will be seen in the near future. Traditional corpus statistics, such as inverse document frequency (*idf*), have different characteristics when documents are encountered one-at-a-time. These issues are important, because they determine how efficient and effective statistical document filtering systems will be in "real world" environments.

10 Document filtering, also known as selective dissemination of information (SDI) is generally based on an unranked Boolean retrieval model. A user's information need is expressed by a query, also called a profile, in a query language. Sometimes a profile is actually a set of queries for one user; in this discussion, query and profile are considered synonymous. Queries are typically expressed for
15 the purposes of this discussion using Boolean logic. A query either matches or does not match a document. There is no ability to partially satisfy a query, or to determine how well a document matches or satisfies a query. Instead, the emphasis is on speed, and on indexing methods that enable very fast processing of documents against profiles.

20 In one example of these classes of systems, each Boolean profile is analyzed to identify the least frequent trigram (LFT) that *must* occur whenever the profile matches a document (a necessary, but not sufficient, condition for matching). Documents are converted into a restricted alphabet, and represented as a sequence of trigrams. For each profile, a table lookup determines whether its
25 LFT is present. If not, the profile can not possibly match the document. This first stage is designed to eliminate greater than 95% of the profiles in just a few

00042757 004897

instructions each. If a profile's LFT is present, a slower second stage determines whether the document actually satisfies the Boolean query.

It is generally accepted that statistical systems provide better precision and recall for document retrieval than do unranked Boolean systems. The growing
 5 power of computer hardware has made statistical systems increasingly practical for even large scale document filtering environments. A common approach has been to simulate document filtering with an existing vector-space or probabilistic document retrieval system on a collection of new or recent documents. This approach is simple, effective, and has the advantage of a corpus from which to
 10 gather statistics like *idf*. However, it is not well-suited to immediate dissemination of new information, and it adds index creation, storage, and maintenance to the cost of document filtering.

Summary of the Invention

A new profile selection technique for structured queries is described which
 15 has a dramatic effect on filtering speed. Prior to filtering, each profile is analyzed to determine the number of document terms it must match before a document can exceed the dissemination threshold. The "optimal" document for a query term is one in which *ntf* approaches 1. The estimate for the minimum number of terms required is made by setting *ntf* to 1 for each query term, and then ordering sibling
 20 query net nodes by the estimated belief values. When reordering is complete, the query net is traversed, accumulating belief values and counting query terms. When the accumulated belief exceeds the dissemination threshold, the minimum number of terms necessary to exceed the threshold is known. This information can be stored in the profile index, and used during document filtering. The profile is
 25 selected only if it matches a sufficient number of document terms.

In one aspect, the present invention relates to an apparatus for filtering documents as those documents are received. The apparatus includes a document parser which accepts incoming documents as input and outputs inverted lists of terms contained in the documents. The apparatus also includes a profile parser
 5 which accepts user queries as input and provides a query net representing user queries as output. A comparator is provided that compares the inverted list for an incoming document against a query net representing a user query and provides an output and indication of whether the incoming document matches the user query.

In another aspect, the present relates to a method for filtering incoming
 10 documents. The method begins by receiving an incoming document and parsing it to produce an inverted list of terms contained in the incoming document. The inverted list is then used to retrieve user queries and user queries matching less than a predetermined number of terms are immediately discarded. Each remaining query is scored and those user queries having a score less than the predetermined
 15 threshold are discarded.

Brief Description of the Drawings

The invention is pointed out with particularity in the appended claims. The advantages of the invention described above, as well as further advantages of the invention, may be better understood by reference to the following description
 20 taken in conjunction with the accompanying drawings, in which:

FIG. 1 is a block diagram of the system of the present invention;

FIG. 2 is a graph showing mean squared error as a function of the number of documents encountered;

FIG. 3 is a graph showing the number of unique word stems encountered in
 25 a corpus as a function of the number of documents encountered;

FIG. 4 is a graph showing average document length as a function of the

number of documents encountered; and

FIG. 5 is a graph showing the average document length as a function of the number of documents encountered.

Detailed Description of the Invention

5 Referring now to FIG. 1, a document filtering system 10 based upon the inference network model of information retrieval and filtering is shown. The major tasks performed by the system are creation of query networks 12, creation of document networks (clipsets) 14, and use of the clipsets 14 to filter documents. The document network is created automatically by mapping documents onto
10 content representation nodes, which are implemented with traditional inverted lists 16. Query networks are specified by a user using either natural language or a structured query language profiles 18. Although the query language may be a traditional Boslean language, it is preferred to use a query language which includes probabilistic AND, OR, and NOT operator, proximity operators, probabilistic
15 phrase and passage operators, and a weighted sum operator for user-specified weights. Document filtering is performed by using recursive inference to propagate belief values through the inference net, discarding any documents whose belief is below a dissemination threshold.

A query network 12 is a directed acyclic graph (DAG) in which the root is
20 a query operator (e.g. #SUM), internal nodes correspond to nested query operators (e.g., AND, OR, phrase, proximity, or other operators), and leaves correspond to query terms. When a supplied user profile 18 is completely parsed, the DAG is optimized, by removing redundant query terms and operators, reordering arguments to Boolean operators, minimize the cost of evaluating the
25 query.

5
10

15

20

25

5

10

15

20

25

discarded. Inverted lists 16 are constructed, incrementally as tokens are encountered, for the remaining stems. When the document is parsed completely, the result is a set of inverted lists representing the document network for that document. Finally, each list is annotated with the belief that the term will

5 contribute. Belief is calculated using *tf.idf* formula, as shown below.

$$n\textit{tf} = 0.04 + 0.6 \cdot \frac{tf}{tf + 0.5 + 1.5 \cdot \frac{dl}{avg_dl}}$$

$$idf = \frac{\log\left(\frac{C + 0.5}{df}\right)}{\log(C + 1.0)}$$

$$bel_{term}(t) = 0.4 + 0.6 \cdot n\textit{tf} \cdot idf$$

where:

- tf* is the frequency of term *t* in the document,
- 10 *dl* is the document length,
- avg_dl* is the average document length in the collection,
- C* is the number of documents in the collection, and
- df* is the number of documents in which term *t* occurs.

Three of the statistics above are derived from the corpus as a whole: *df*, *avg_dl*,
 15 and *C*. Accurate values for these three statistics are known only after all documents are filtered, so filtering generally must be performed only with estimates.

Parsing a 3,000 byte document and later freeing the associated indices and data structures using the described method takes 0.02937 seconds of wall-clock
 20 time on an otherwise idle DECStation 3000-6000 equipped with 64 megabytes of memory and running at 175 Mhz, manufactured by Digital Equipment Corporation of Maynard, Massachusetts.

Document parsing speed is affected by the number of profiles 18, because inverted lists 16 are built only for terms in the profile term dictionary 32. As more profiles 18 are added, the vocabulary grows larger. Fortunately, adding a large number of profiles 18 causes only a small increase in the size of the term dictionary 32, and therefore only a small decrease in document parsing speed.

After a document is indexed, it can be compared to a clipset 40. Retrospective document retrieval systems owe their speed partially to indexing methods, such as inverted lists, that enable the system to consider only those documents that have terms in common with a query. A similar need exists for document filtering, because many profiles have nothing in common with most documents.

Once a set of profiles is selected, each profile must be compared to an incoming document. One embodiment iterates through the selected profiles, determining for each the belief that the document satisfies the information need. The belief in a document for a particular query net is determined with depth-first evaluation. For each query term that occurs in the document, this embodiment must locate the appropriate inverted list, lookup the belief associated with the term, and then combine the belief with the beliefs from other terms, according to the probabilistic operators being used. If proximity operators are used, this embodiment must also lookup the locations where the term occurs, intersect those with the locations of other proximate terms, and then compute the belief for the proximity operator dynamically, using the same *tf.idf* formulas described above.

A profile is returned for a document if and only if it matches the query, *and* if the belief that the document satisfies the information need exceeds a profile-specific, user-supplied dissemination threshold. This latter requirement is particularly important for probabilistic query operators. A document “matches” a

weighted sum or probabilistic (“fuzzy”) AND operator if even one query term is present, although the belief in the document is usually low. This behavior is rarely a problem in a ranked retrieval system with large sets of documents, because a low belief causes a document to appear low in the rankings. However, in a filtering
5 environment, a low-scoring document may still be the best document encountered that day. If the system does not discard documents with low beliefs, users must either develop strictly Boolean queries, or wade through irrelevant documents on days when no relevant documents occur.

Filtering a 3,000 byte document for 1,000 Boolean profiles, each
10 containing an average of 22 terms and operators, takes 0.024 seconds of wall-clock time on an otherwise idle DECStation 3000-600, manufactured by Digital Equipment Corporation of Maynard, Massachusetts. The time is proportional to the number of profiles; twice as many profiles takes twice as long. In one experiment a 109 megabyte file of 39,906 Wall Street Journal documents was
15 processed against 1,000 Boolean profiles in 25 minutes of wall-clock time, generating 616,487 matches. A similar experiment with 1,000 statistical profiles of similar complexity required 33 minutes.

EXPERIMENTS

20 One important difference between document filtering and document retrieval is how corpus-wide statistics like inverse document frequency (*idf*) and average document length are obtained. The effectiveness of current retrieval models depends upon accurate corpus statistics, which document retrieval systems gather while indexing the collection. In an on-line environment, where documents
25 must be filtered as soon as they arrive, accurate corpus statistics are not available until after all of the documents have been filtered.

One approach used in some experiments was to use statistics from another, presumably similar, corpus. This approach is effective, but it may be impractical in practice. Obtaining *idf* values from a retrospective corpus can be expensive, particularly if queries include large numbers of proximity operators. *Idf* values for
 5 unindexed query fragments (e.g. proximity operators) can only be obtained by running queries against the retrospective collection. In one experiment, it took several hours to obtain the *idf* values for proximity operators in 50 routing queries. This cost would be prohibitive in a “real world” setting.

An alternate approach is to estimate corpus statistics dynamically as each
 10 document is encountered. This approach has the advantage of being “low cost” and of not requiring a similar training corpus. Although the corpus statistics will initially be inaccurate, they eventually converge to their “true” values for the corpus.

Experiments were performed to study *idf* values on a small corpus, because
 15 it converges relatively quickly. Average document length was studied on a larger corpus, because it converges less quickly.

Figure 2 shows the convergence of *idf* values for terms in the 1988 Wall Street Journal corpus. Each curve shows the mean square error (MSE) between estimated *idf* and true *idf* at 1,000 document increments. The top curve shows the
 20 MSE for a traditional method of computing *idf*. The bottom curve shows the MSE for the “scaled” *idf* used by the INQUERY document retrieval product manufactured by Sovereign Hill Software, Inc. of Dedham, Massachusetts.

Idf values converged rapidly to their true values (Figure 2), even as the vocabulary continued to grow (Figure 3). A “scaled” *idf* converges more rapidly,
 25 because it gives a more accurate estimate for terms that occur just once. An unscaled *idf* for terms that occur just once changes significantly as more

documents are observed, while a scaled *idf* changes very little. If terms that occur just once are excluded (middle curve), the MSE for the traditional method is reduced by about half.

Figure 4 shows the convergence of average document length for the TREC-4 Routing corpus. It takes about 25,000 documents to reach a stable estimate, but the estimate then changes significantly whenever the document stream shifts from one subcollection to another.

The effect of shifting from one subcollection to another can be eliminated by interleaving the subcollections. Although the documents could be ordered by publication date, doing so does not eliminate the “subcollection” effect because the subcollections cover different periods of time. Figure 5 shows the convergence of average document length in a proportionally interleaved TREC-4 Routing corpus. It takes about 20,000 documents to reach a stable estimate in this corpus, but the estimate is 15% above its eventual final value, and it continues to drift up and down, smoothly but by significant amounts, for another 100,000 documents.

Precision	Number of Documents Used Only For Training						
	0	1000	3000	5000	10000	15000	20000
at 5 docs	-18.7%	-16.1%	-11.8%	-8.1%	-3.1%	-3.7%	-3.1%
at 10 docs	-16.4%	-13.8%	-10.9%	-9.0%	-3.9%	-2.9%	-3.5%
at 15 docs	-16.0%	-11.9%	-9.9%	-7.0%	-3.5%	-2.2%	-1.1%
at 20 docs	-14.7%	-13.2%	-10.5%	-8.5%	-5.4%	-3.8%	-2.7%
at 30 docs	-13.6%	-12.2%	-10.0%	-9.1%	-5.7%	-3.5%	-3.4%
at 100 docs	-9.9%	-11.2%	-9.2%	-8.1%	-5.8%	-5.3%	-4.7%
at 200 docs	-2.8%	-9.1%	-8.1%	-7.3%	-5.3%	-5.3%	-5.3%
at 500 docs	+4.4%	-3.9%	-3.3%	-3.0%	-2.8%	-4.0%	-5.3%
Avg	-8.0%	-12.1%	-10.7%	-10.7%	-6.9%	-7.2%	-7.8%

TABLE 1

An experiment with TREC-4 Routing queries and documents investigated the effects on recall and precision of learning corpus-wide statistics during filtering.

InFilter, a routing component of the INQUERY product, was run twice on the TREC-4 corpus (935 MB, 329,780 documents) and INQ203 Routing queries (50 queries, 50 terms and 200 proximity pairs each). In one run, corpus statistics were available *a priori* (“perfect” statistics). In the other, estimates were updated as each document was encountered (“learned” statistics). The experiment required dissemination thresholds that would disseminate at least 1,000 documents for each query. We used the documents score that INQUERY assigned at rank 1,000, because it was conveniently available.

Learned corpus statistics produced a significant loss in average precision at all cutoffs and levels of recall (Table 1, Column “0”). The effect of inaccurate corpus statistics in the first few thousand documents is rather dramatic, given that the estimates converge to relatively accurate values after filtering only a small percentage of the corpus. However, analysis reveals that learned statistics produce substantially higher scores for documents filtered “early” than for documents filtered “later”, when corpus statistics have converged. The “early” documents, with their overly generous scores, dominate the top of the rankings.

If the first several thousand documents are used only for training purposes (i.e., are not disseminated), the effect of learned corpus statistics on recall and precision is less significant (Table 1, columns “1000” to “20000”). For example, if 15,000 documents are used for training, corpus statistics produce a 2.5-5.3% loss in precision at cutoffs 5-500. This is a crude way of analyzing the effects of learning corpus statistics, because the baseline is based on all of the relevant documents, while the filtered set is missing whatever relevant documents were

discarded during training. However, it confirms that, after the initial period of training, learned corpus statistics are effective for filtering.

Speed is an important characteristic of document filtering systems, and consequently techniques for optimizing Boolean filtering systems are well-known.

5 Similar techniques for statistical document filtering are required.

Filtering a document involves profile selection and evaluation. Profile selection determines which profiles to evaluate; profile evaluation determines how well a document satisfies a profile. For each document, the goal is to spend either no time or nearly no time on most of the profiles.

10 One approach is to index profiles with inverted lists. The terms in a document “retrieve” profiles during filtering. This approach works particularly well with the unstructured queries that characterize vector-space systems, because profile scores can be computed when inverted lists are merged.

Profile indexing is less effective with the structured queries that
15 characterize inference network systems, because scores for structured queries cannot be computed when profile inverted lists are merged, that is, the inference network can simply be turned “upside down.” In this case, profile indexing can be used only to identify profiles that are candidates for evaluation. Profile indexing may also be less effective on long routing queries, because a profile with many
20 terms is more likely to have at least one in common with any document.

Using the method described above, an estimate for the minimum number of necessary terms is obtained with an algorithm similar to algorithms that reorder and/or optimize unstructured queries and Boolean queries. Reordering by optimal belief is perhaps a more general technique, because it applies to both unstructured
25 queries and queries structured with a wide range of Boolean and probabilistic operators. However the important difference is that the query is not reordered to

optimize query evaluation (although doing so is a good idea), but to find the minimum number of terms necessary to select a profile.

The method can be described as implemented as a three stage filter. First, document terms “retrieve” profiles, using inverted lists. The number of terms
 5 matching each profile is determined as inverted lists are merged. Next, “retrieved” profiles that don’t match enough document terms are discarded. Finally, the remaining profiles are evaluated completely, and any with scores below the dissemination threshold are discarded.

The speedup obtained with the described method increases as a profile’s
 10 dissemination threshold increases. If the threshold is low, the minimum number of terms necessary to select a profile is one, reducing the method to simple profile indexing.

The described method can be a “safe” or “unsafe” optimization, depending upon how it is used. If profiles are reanalyzed each time the *idf* values change, it is
 15 safe, i.e., guaranteed to select for a given document every profile that can possibly exceed the dissemination threshold. If *idf* values change, as when they are being learned, the estimate may become wrong. Usually the estimate will be an underestimate, causing no harm, because *idf* values can fall rapidly (increasing the actual number of necessary terms), but tend to rise slowly (decreasing the actual
 20 number of necessary terms). However, it may make sense to reanalyze profiles periodically, for example every few thousand documents, when *idf* values are being learned.

The relative effectiveness of these techniques is demonstrated in two experiments. In on experiment, the TREC 1988 Wall Street Journal corpus (109
 25 MB, 39,906 documents) was filtered for a set of 3,000 simple, artificially-generated profiles (10 terms and 4 proximity pairs each). The dissemination

threshold was set to yield about a 0.2% “hit” rate. In the second experiment, the TREC-4 Routing corpus (935 MB, 329,780 documents) was filtered for a set of 50 complex profiles (50 terms and 200 proximity pairs each). The dissemination threshold was set to yield about 1,000 documents per profile (a 0.3% “hit” rate), as is common in TREC Routing evaluations. In both experiments InFilter was learning corpus statistics, so profiles were re-analyzed and their MinTerm estimates updated every 1,000 documents.

With simple profiles (the 1988 WSJ experiment), simple profile indexing was a substantial improvement over evaluating all profiles. Filtering time was reduced by 37.5% without impacting effectiveness. Estimates for the number of necessary terms were updated every 1,000 documents. Table 2 below summarizes the results.

	3,000 simple profiles 1988 WSJ corpus			50 complex profiles TREC-4 Routing corpus		
	No. Index	Inverted Index	Min Term Index	No. Index	Inverted Index	Min Term Index
Profiles Fully Evaluated	100%	24.2%	4.25%	100%	97.5%	74.9%
Total Filtering Time (h:mm)	1:28	0:55	0:41	5:26	5:31	2:53
Filtering Rate (MB / hour)	74	119	160	172	170	324
Avg Documents Disseminated Per Profile	77.8	77.8	77.8	999.5	999.5	922.0

With simple profiles (the 1988 WSJ experiment), simple profile indexing was a substantial improvement over evaluating all profiles. Filtering time was reduced by 37.5% without impacting effectiveness. MinTerm Indexing reduced filtering time by another 25%.

With complex profiles (the TREC-4 Routing experiment), simple profile indexing was slightly *worse* than evaluating all profiles. The computational cost of simple profile indexing provided little benefit, because most documents had a term in common with most of these routing profiles. However MinTerm Indexing,

which considers the *number* of terms a document has in common with a profile, reduced filtering time by 47%.

MinTerm indexing was “unsafe” in these experiments, because corpus statistics were updated after each document but profile MinTerm estimates were
5 updated after each 1,000 documents. In the 1988 WSJ experiment, the cost was the lost of one document from a set of 233,735. In the TREC-4 corpus, the cost was a much higher 76.5 documents per profile.

Most of the TREC-4 loss was due to experimental error. The algorithm that determined the number of terms a document and profile have in common did
10 not consider the effect of duplicate terms in the profile. Duplicates are very rare in the 1988 WSJ profiles, so this error had no effect on the first experiment. Duplicates are quite common in the TREC-4 profiles, hence the “missed” documents in the second experiment.

15 Having described certain embodiments of the invention, it will now become apparent to one of ordinary skill in the art that other embodiments incorporating the concepts of the invention may be used. Therefore, the invention should not be limited to certain embodiments, but rather should be limited only by the spirit and scope of the following claims.

0894357-034997

CLAIMS

What is claimed is:

- 1 1. An apparatus for filtering documents as received, the apparatus
2 comprising:
3 a document parser, said document parser accepting documents as input and
4 providing inverted lists of terms contained in the documents as output;
5 a profile parser, the profile parser accepting user queries as input as
6 providing a query net representing the query as output; and
7 a comparator that compares the inverted list for an incoming document
8 against the query net representing the user query and providing as output an
9 indication whether the incoming document matches the user query.
- 1 2. The apparatus of claim 1 wherein said profile parser provides as output a
2 term dictionary containing terms from the user query, and wherein the document
3 parser uses said term dictionary to eliminate terms from said inverted list.
- 1 3. The apparatus of claim 1 wherein said profile parser accepts a plurality of
2 user queries and stores a corresponding plurality of query nets in a memory
3 element.
- 1 4. The apparatus of claim 3 wherein said comparator compares an inverted
2 list associated with an incoming document against each of said plurality of stored
3 query nets.
- 1 5. The apparatus of claim 1 wherein said comparator provides an indication
2 that the incoming document matches the user query if said inverted lists contain a
3 minimum number of terms contained in the query net.

0894257.081897

1 6. The apparatus of claim 1 wherein said profile parser, said document parser,
2 and said comparator reside on separate machines, said separate machines
3 interconnected by a network.

1 7. A method for filtering incoming documents, the method comprising the
2 steps of:

3 (a) receiving an incoming document and parsing it to produce an
4 inverted list of terms contained in incoming document;

5 (b) using the produced inverted list to retrieve query nets representing
6 user queries;

7 (c) discarding retrieved query nets matching less than a predetermined
8 number of terms;

9 (d) scoring remaining profile and discarding profiles having a score less
10 than a predetermined threshold.

1 8. The method of claim 7 further comprising the step of receiving user queries
2 and parsing the user queries to produce query nets representing the query.

1 9. The method of claim 7 further comprising that step of storing all query nets
2 associated with a user as a clipset.

1 10. The method of claim 7 further comprising the step of storing an inverted
2 list associated with an incoming document.

1 11. The method of claim 9 further comprising the step of storing corpus
2 statistics in the clipset.

1 12. An article of manufacture having computer-readable program means for
2 filtering incoming documents, the article comprising:

03942767 "031697"
469789 4524680

- 3 (a) computer-readable program means for receiving an incoming
4 document and parsing it to produce an inverted list of terms contained in incoming
5 document;
- 6 (b) computer-readable program means for using the produced inverted
7 list to retrieve query nets representing user queries;
- 8 (c) computer-readable program means for discarding retrieved query
9 nets matching less than a predetermined number of terms;
- 10 (d) computer-readable program means for scoring remaining profile
11 and discarding profiles having a score less than a predetermined threshold.

1 13. The article of manufacture of claim 12 further comprising computer-
2 readable program means for receiving user queries and parsing the user queries to
3 produce query nets representing the query.

1 14. The article of manufacture claim 12 further comprising computer-readable
2 means of storing all query nets associated with a user as a clipset.

1 15. The article of manufacture of claim 12 further comprising
2 computer-readable means of storing an inverted list associated with an incoming
3 document.

1 16. The article of manufacture of claim 9 further comprising computer-readable
2 means of storing corpus statistics in the clipset.

SYSTEM AND METHOD FOR FILTERING A DOCUMENT STREAM

Abstract of the Disclosure

A system for filtering documents and includes a document parser, a profile parser, and a comparator. The document parser accepts incoming documents as input and provides inverted lists of terms contained in the document's output. The profile parser accepts as input user queries and provides as output query nets representing the user queries. The comparator compares the inverted lists representing the documents against the query that is representing the user queries to determine if an incoming document matches a user query. A related method for filtering incoming documents includes the steps of receiving an incoming document and parsing it to produce an inverted list of terms contained in the incoming document. The inverted list is then used to retrieve user queries. Any user queries matching less than a pre-determined number of terms are immediately discarded. The remaining user queries are scored and user queries having a score less than a predetermined threshold are discarded. The remaining user queries are the queries which the incoming document matches.

15

20

**COMBINED DECLARATION AND POWER OF ATTORNEY
FOR PATENT APPLICATION**

(Original, Design, National Stage of PCT, Supplemental, Divisional, Continuation or CIP)

As a below named inventor, I hereby declare that:

My residence, post office address and citizenship are as stated below next to my name, and I believe I am the original, first and sole inventor (if only one name is listed below) or an original, first and joint inventor (if plural names are listed below) of the subject matter which is claimed and for which a patent is sought on the invention entitled:

SYSTEM AND METHOD FOR FILTERING A DOCUMENT STREAM

the specification of which (check one):

- ☒ is attached hereto.
- ☐ was filed on _____ and identified by Attorney Docket No. SVH-003 (7453/6) and accorded Application Serial No. 0 / or
- ☐ was described and claimed in PCT International Application No. _____ filed on _____ and as amended under PCT Article 19 on _____ (if any).

I hereby state that I have reviewed and understand the contents of the above identified specification, including the claims as amended by any amendment referred to herein.

I acknowledge the continuing duty to disclose information which is material to the examination of this application in accordance with 37 C.F.R. §1.56.

PRIORITY CLAIM

- ☐ A. I hereby claim benefit under 35 U.S.C. 119(e) of United States Provisional Application No. _____, filed on _____.
- ☐ B. I hereby claim foreign priority benefits under 35 U.S.C. §119 of any foreign application(s) for patent or inventor's certificate or of any PCT international application(s) designating at least one country other than the United States of America listed below and I have also identified below any foreign application(s) for patent or inventor's certificate or any PCT international application(s) designating at least one country other than the United States of America filed by me on the same subject matter having a filing date before that of the application(s) of which priority is claimed.

Express Mail Label No. EM353713966US

- ☐ no such applications have been filed.
- ☐ such applications have been filed as follows:

**EARLIEST FOREIGN APPLICATION(S), IF ANY FILED WITHIN
12 MONTHS (6 MONTHS FOR DESIGN) PRIOR TO
THIS U.S. APPLICATION**

Country	Application Number	Date of Filing (mo., day, year)	Priority Claimed Under 35 USC 119
			<input type="checkbox"/> YES NO <input type="checkbox"/>
			<input type="checkbox"/> YES NO <input type="checkbox"/>
			<input type="checkbox"/> YES NO <input type="checkbox"/>

- ☐ C. I hereby claim the benefit under 35 U.S.C. §120 of any United States application(s) or PCT international application(s) designating the United States of America that is/are listed below and, insofar as the subject matter of each of the claims of this application is not disclosed in that/those prior application(s) in the manner provided by the first paragraph of 35 U.S.C. §112, I acknowledge the duty to disclose material information as defined in 37 C.F.R. §1.56 which occurred between the filing date of the prior application(s) and the national or PCT international filing date of this application.

**PRIOR U.S. NON-PROVISIONAL APPLICATIONS OR PCT INTERNATIONAL
APPLICATIONS DESIGNATING THE U.S. FOR BENEFIT UNDER 35 USC §120:**

U.S. APPLICATIONS	U.S. FILING DATE	STATUS
(Application Serial No.)	(Filing Date)	(Status) (patented, pending, aband.)
(Application Serial No.)	(Filing Date)	(Status) (patented, pending, aband.)
(Application Serial No.)	(Filing Date)	(Status) (patented, pending, aband.)

POWER OF ATTORNEY

As a named inventor, I hereby appoint the following attorneys and/or agents to prosecute this application and transact all business in the United States Patent and Trademark Office connected therewith:

Steven M. Bauer	Reg. No. 31,481
Paula A. Campbell	Reg. No. 32,503
Joseph A. Capraro, Jr.	Reg. No. 36,471
John J. Cotter	Reg. No. 38,116
Gillian M. Fenton	Reg. No. 36,508
Duncan A. Greenhalgh	Reg. No. 38,678
William G. Guerin	Reg. No. P41,047
Douglas J. Kline	Reg. No. 35,574
John D. Lanza	Reg. No. 40,060
Timothy P. Linkkila	Reg. No. P-40,702
Robin R. Longo	Reg. No. 40,071
Thomas C. Meyers	Reg. No. 36,989
Edmund R. Pitcher	Reg. No. 27,829
Kurt Rauschenbach	Reg. No. 40,137
Michelle B. Rosenberg	Reg. No. P-40,792
J. Scott Southworth	Reg. No. 39,382
Christopher W. Stamos	Reg. No. 35,370
Thomas A. Turano	Reg. No. 35,722
Michael J. Twomey	Reg. No. 38,349
Christine C. Vito	Reg. No. 39,061
Patrick R.H. Waller	Reg. No. P-41,418

Direct correspondence to:

Patent Administrator
Testa, Hurwitz & Thibault, LLP
High Street Tower
125 High Street
Boston, MA 02110

Direct telephone calls to:

John D. Lanza (617) 248-7604

I hereby declare that all statements made herein of my own knowledge are true and that all statements made on information and belief are believed to be true; and further that these statements were made with the knowledge that willful false statements and the like so made are punishable by fine or imprisonment, or both, under Section 1001 of Title 18 of the United States Code, and that such willful false statements may jeopardize the validity of the application or any patent issued thereon.

SIGNATURES

Jamie Callan

Full name of inventor

USA

Citizenship

Inventor's signature

Date

Residence

Same as above

Post Office Address

412JDL7453/6.398533-1

089437034397
463730 46221680

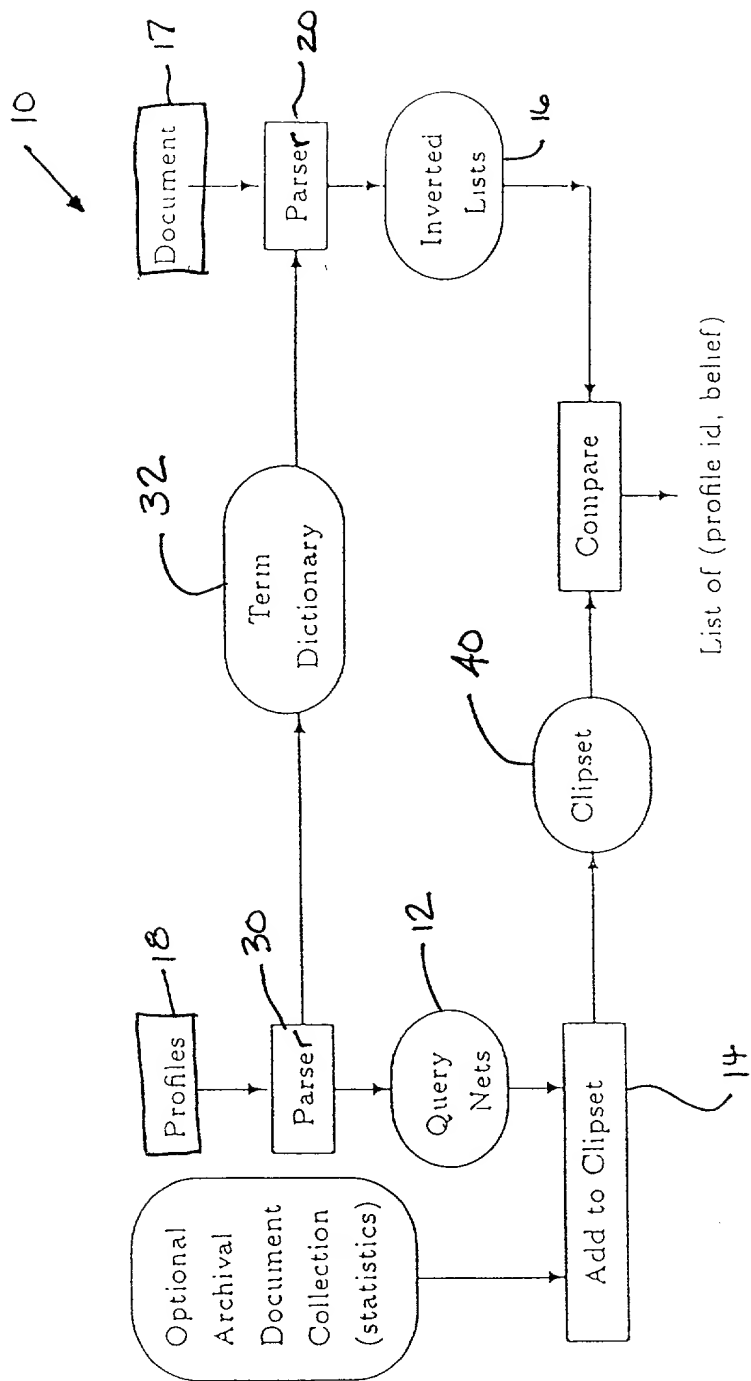


FIG. 1

